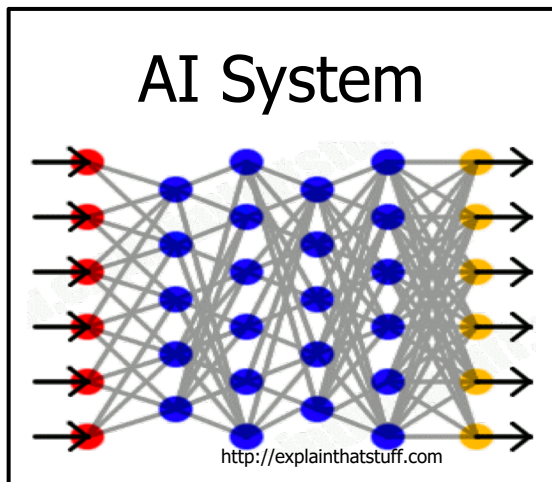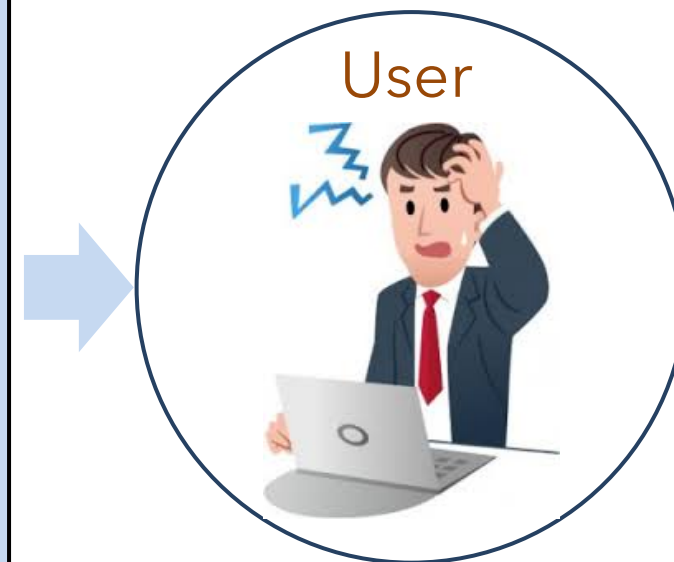# Explainable Artificial Intelligence (XAI)
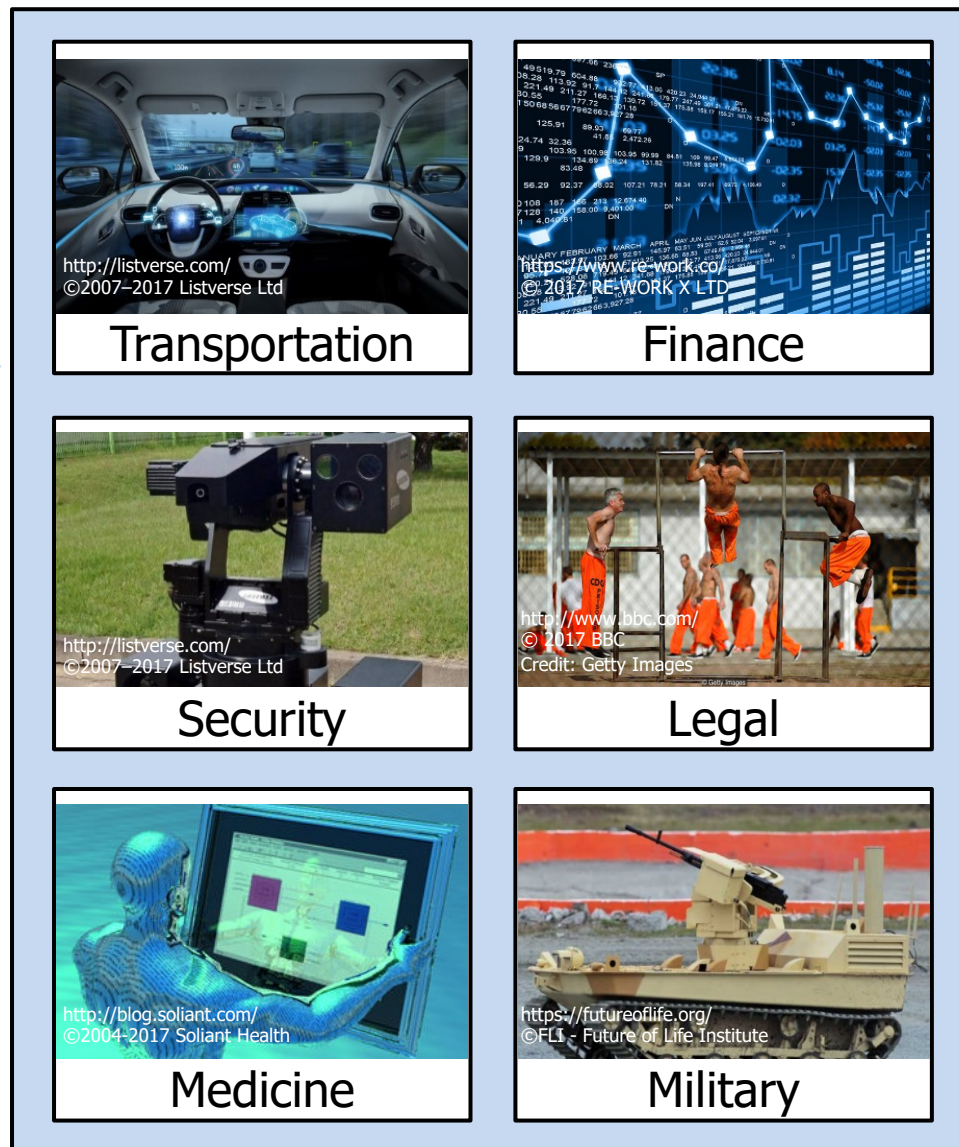


David Gunning
Information Innovation Office (I2O)
Defense Advanced Research Projects Agency (DARPA)

# The Need for Explainable AI

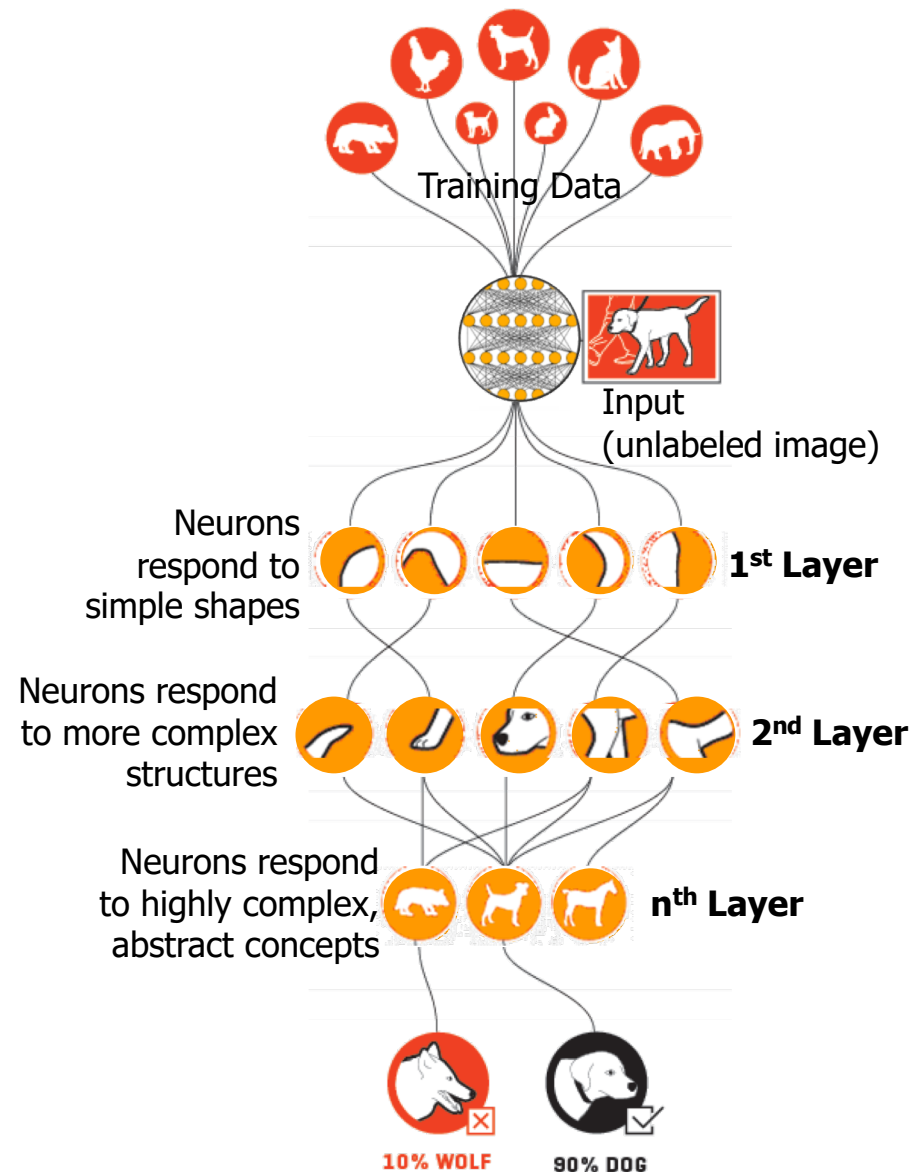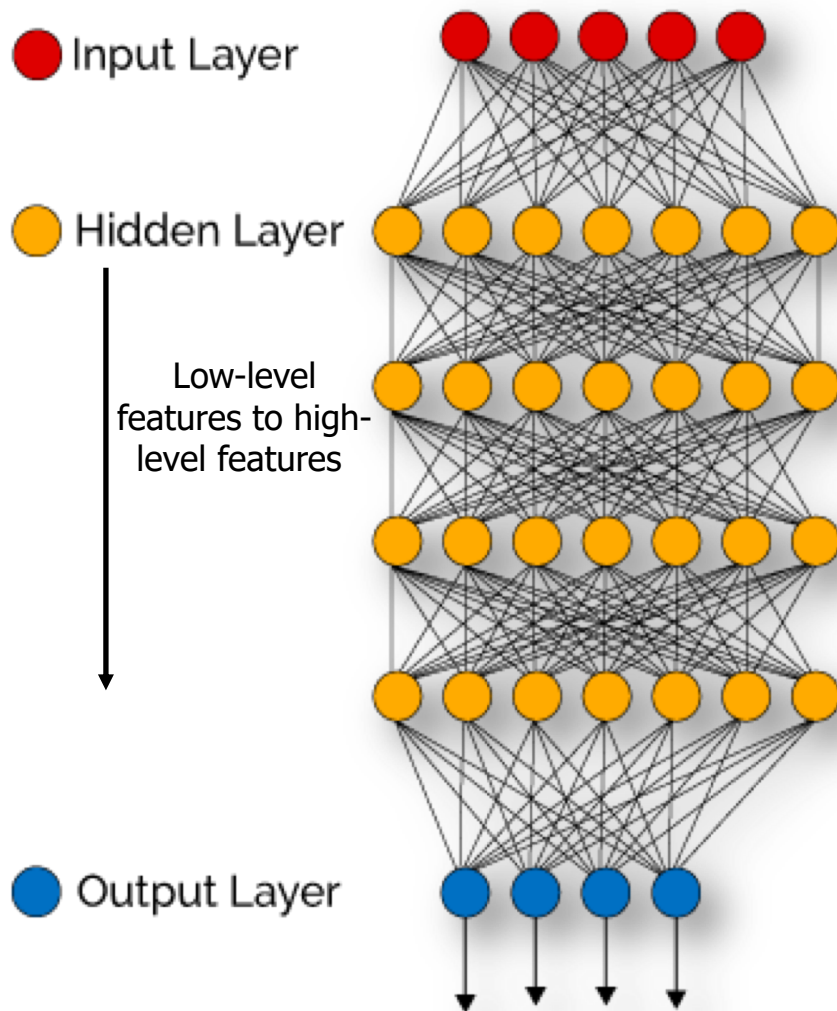## AI System


http://explainthatstuff.com

- We are entering a new age of AI applications
- Machine learning is the core technology
- Machine learning models are opaque, non-intuitive, and difficult for people to understand


http://listverse.com/
©2007–2017 Listverse Ltd
**Transportation**


https://www.re-work.co/
© 2017 RE-WORK X LTD
**Finance**


http://listverse.com/
©2007–2017 Listverse Ltd
**Security**


http://www.bbc.com/
© 2017 BBC
Credit: Getty Images
**Legal**


http://blog.soliant.com/
©2004-2017 Soliant Health
**Medicine**


https://futureoflife.org/
©FLI - Future of Life Institute
**Military**

## User



- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

Input Layer

Hidden Layer

Low-level features to high-level features

Output Layer

Training Data

Input (unlabeled image)

Neurons respond to simple shapes — **1st Layer**

Neurons respond to more complex structures — **2nd Layer**

Neurons respond to highly complex, abstract concepts — **nth Layer**

**10% WOLF**

**90% DOG**

# What are we trying to do?

## Today



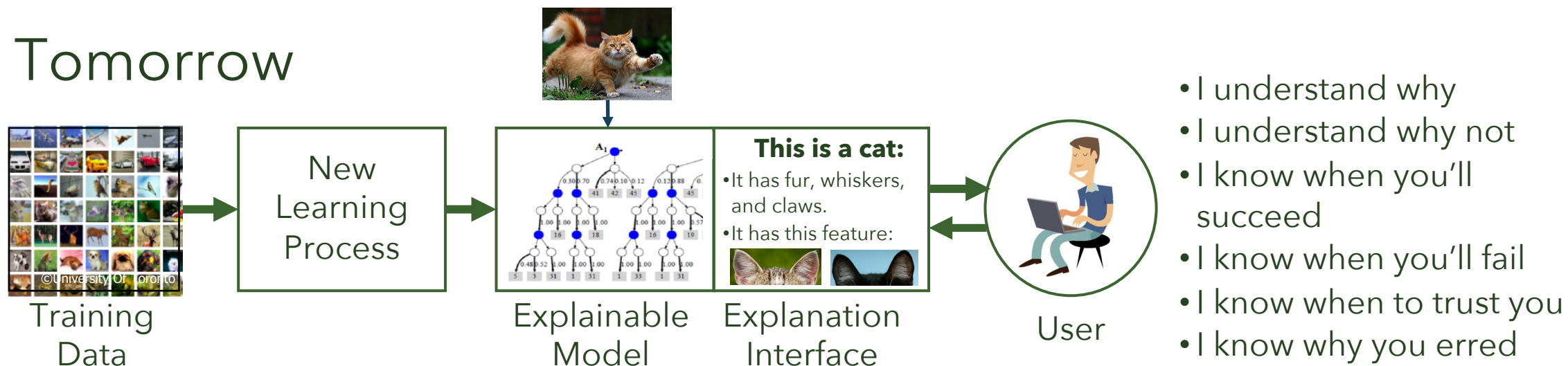Training Data → Learning Process → Learned Function → Output: **This is a cat** (p = .93) → User

- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

## Tomorrow



Training Data → New Learning Process → Explainable Model → Explanation Interface: **This is a cat:**
- It has fur, whiskers, and claws.
- It has this feature:

→ User

- I understand why
- I understand why not
- I know when you'll succeed
- I know when you'll fail
- I know when to trust you
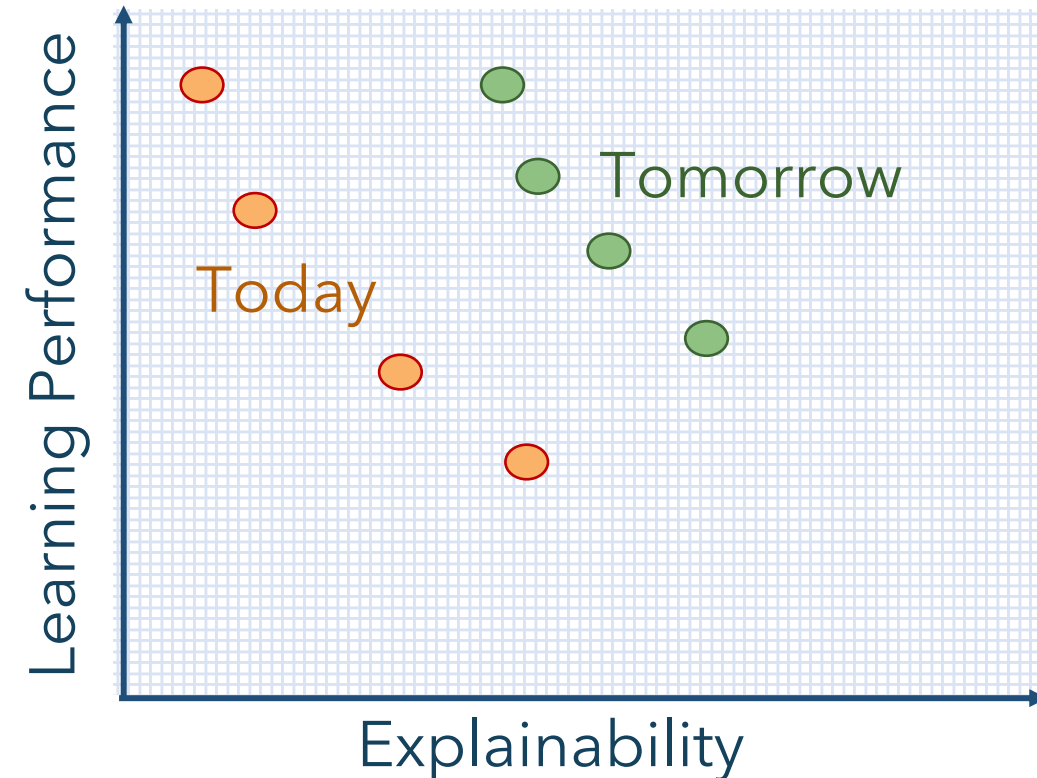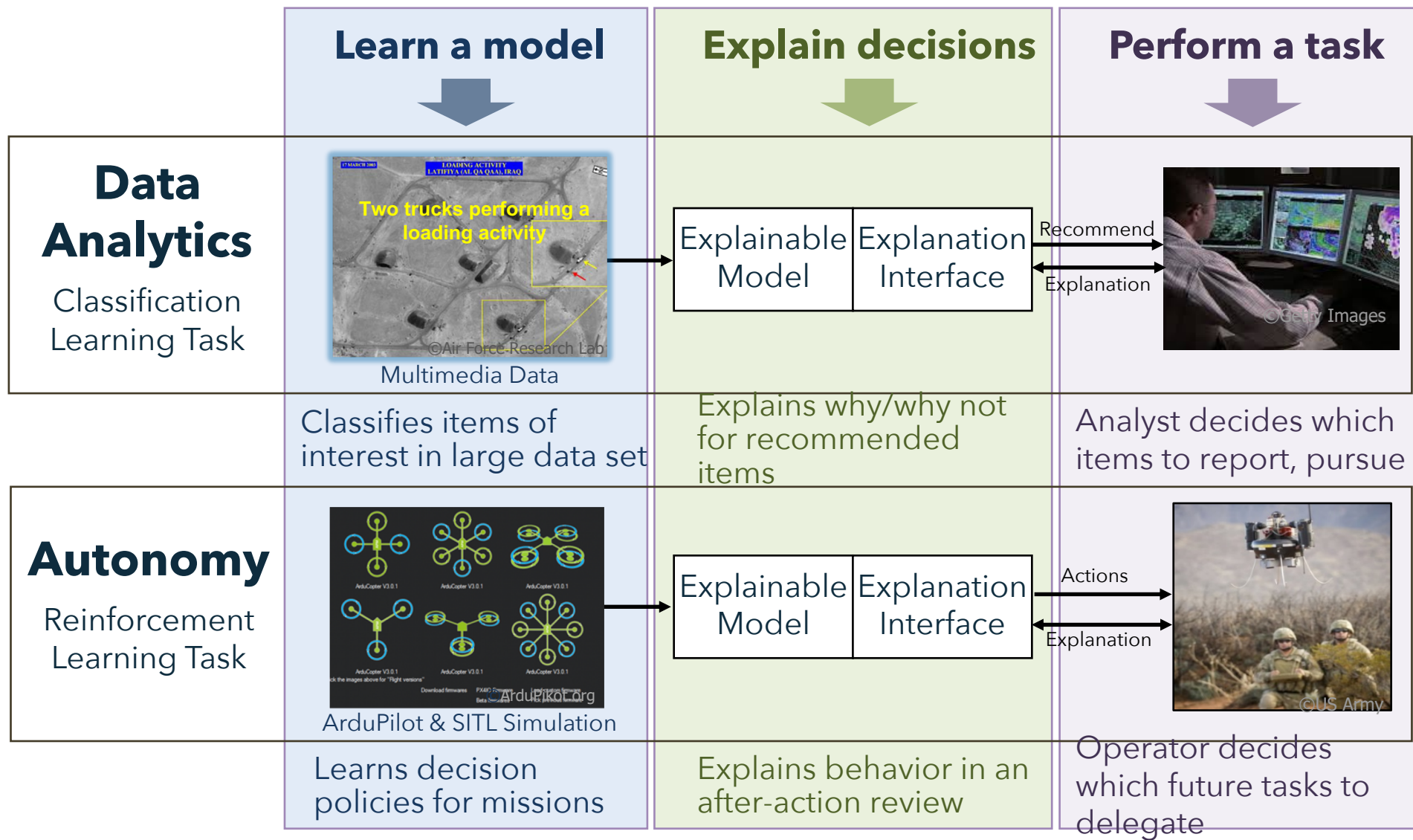- I know why you erred

# Goal: Performance and Explainability

- XAI will create a suite of machine learning techniques that
  - Produce more explainable models, while maintaining a high level of learning performance
  - Enable human users to understand, appropriately trust, and effectively manage the emerging generation of AI systems
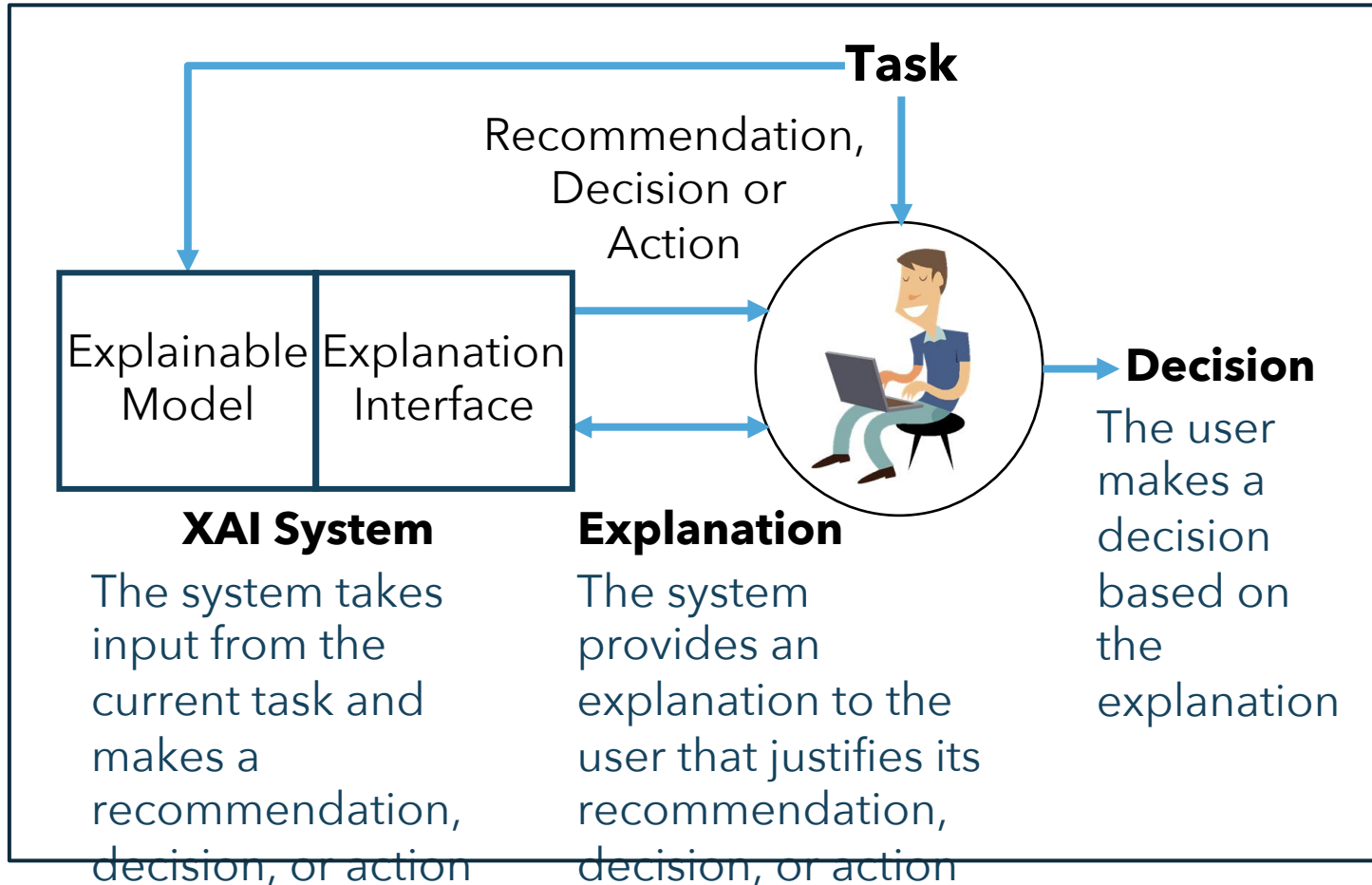
# Challenge Problems



| | **Learn a model** | **Explain decisions** | **Perform a task** |
|---|---|---|---|
| **Data Analytics**<br>Classification Learning Task | Multimedia Data | Explainable Model / Explanation Interface → Recommend / ← Explanation | An analyst is looking for items of interest in massive multimedia data sets |
| | Classifies items of interest in large data set | Explains why/why not for recommended items | Analyst decides which items to report, pursue |
| **Autonomy**<br>Reinforcement Learning Task | ArduPilot & SITL Simulation | Explainable Model / Explanation Interface → Actions / ← Explanation | An operator is directing autonomous systems to accomplish a series of missions |
| | Learns decision policies for missions | Explains behavior in an after-action review | Operator decides which future tasks to delegate |

# Measuring Explanation Effectiveness

## Explanation Framework

**Task**

Recommendation, Decision or Action

**XAI System**
Explainable Model | Explanation Interface

The system takes input from the current task and makes a recommendation, decision, or action

**Explanation**
The system provides an explanation to the user that justifies its recommendation, decision, or action

**Decision**
The user makes a decision based on the explanation

### User Satisfaction
- Clarity of the explanation (user rating)
- Utility of the explanation (user rating)

### Mental Model
- Understanding individual decisions
- Understanding the overall model
- Strength/weakness assessment
- 'What will it do' prediction
- 'How do I intervene' prediction

### Task Performance
- Does the explanation improve the user's decision, task performance?

### Trust Assessment
- Appropriate future use and trust

### Correctability (Extra Credit)
- Identifying errors
- Correcting errors

# Developing an Explainable Model



Learning Techniques (today)

Neural Nets
Deep Learning
Graphical Models
Ensemble Methods
Bayesian Belief Nets
SRL
Random Forests
CRFs      HBNs
MLNs
Statistical Models
AOGs
Decision Trees
SVMs
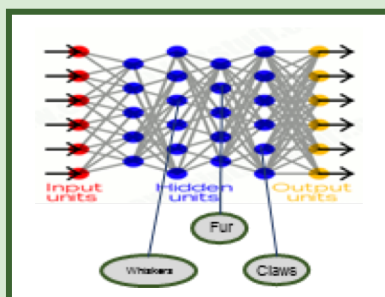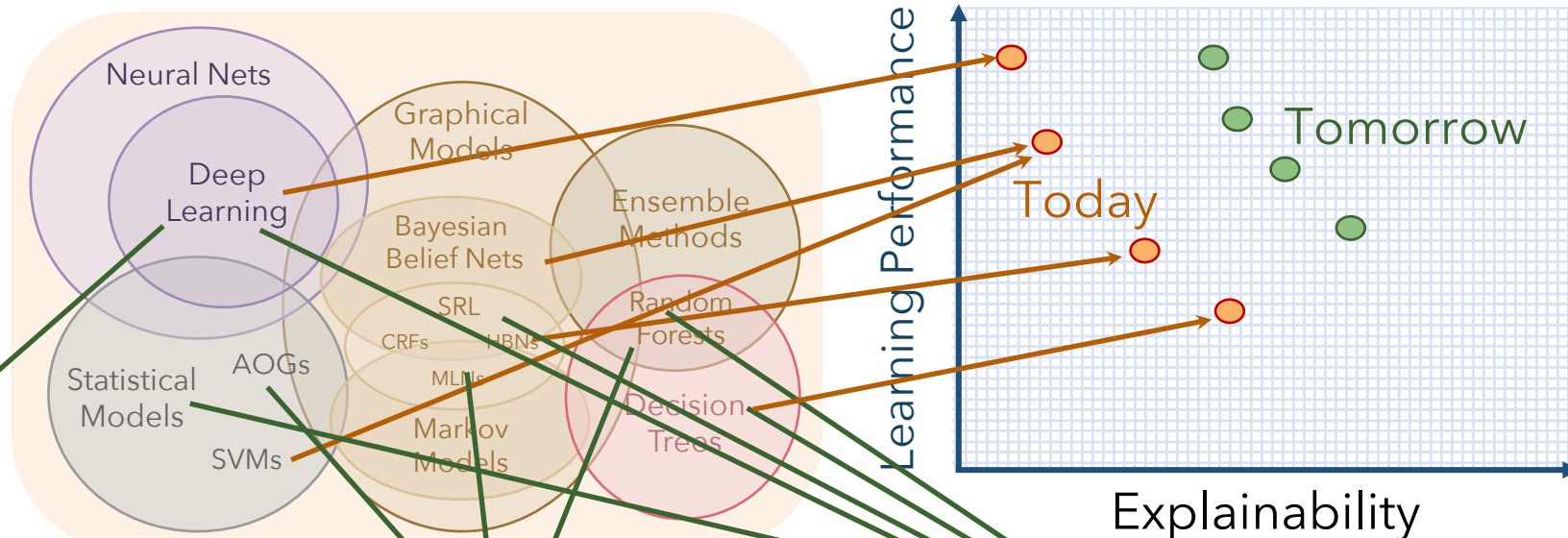Markov Models

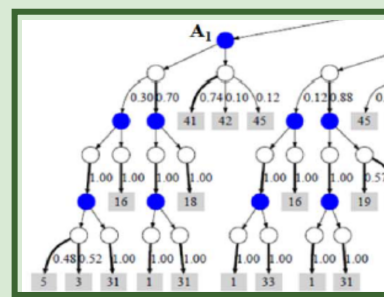Learning Performance

Today

Explainability

Learning Techniques (today)

**XAI Goal**

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance
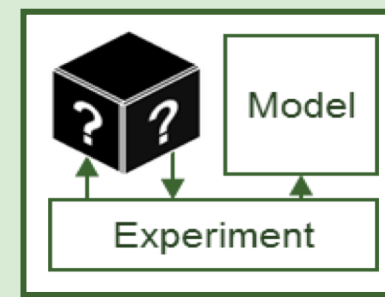
Neural Nets
Graphical Models
Deep Learning
Ensemble Methods
Bayesian Belief Nets
SRL
Random Forests
CRFs    HBNs
Statistical Models    AOGs
MLNs
Decision Trees
SVMs
Markov Models

Learning Performance

Tomorrow

Today

Explainability

**Deep Explanation**
Modified deep learning techniques to learn explainable features
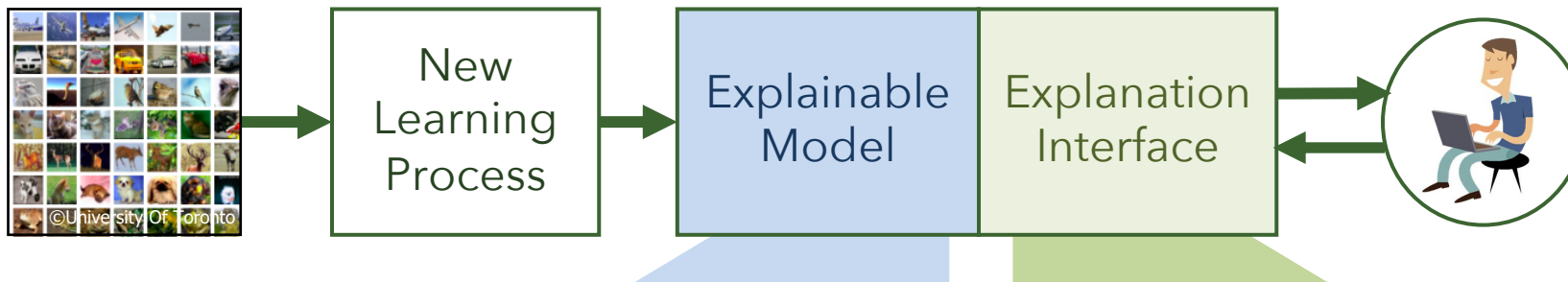
**Interpretable Models**
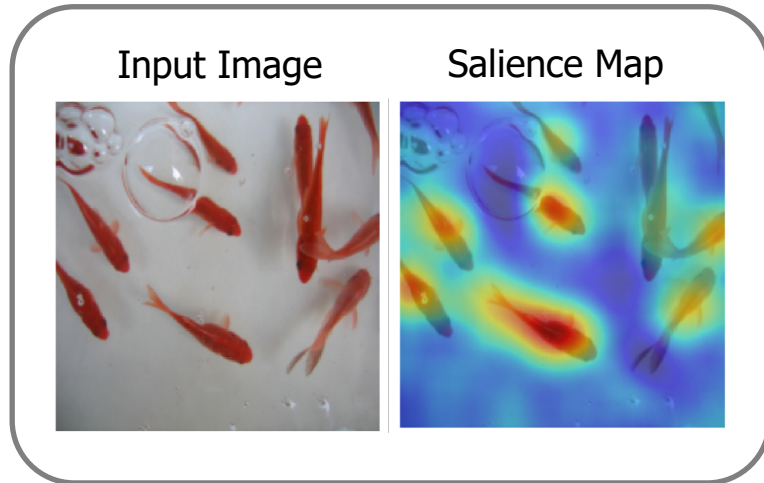Techniques to learn more structured, interpretable, causal models

**Model Induction**
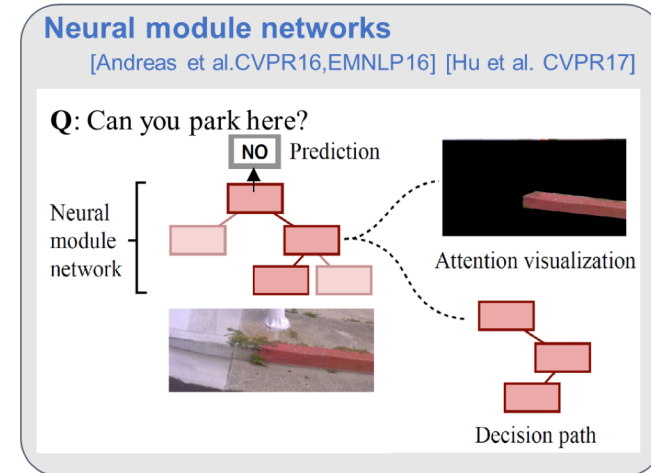Techniques to infer an explainable model from any model as a black box

# XAI Developers and Technical Approaches



IHMC
Psychological Models
of Explanation

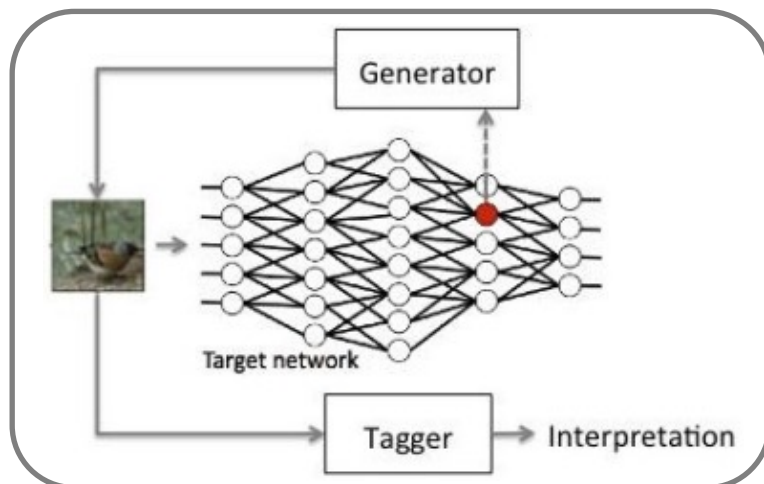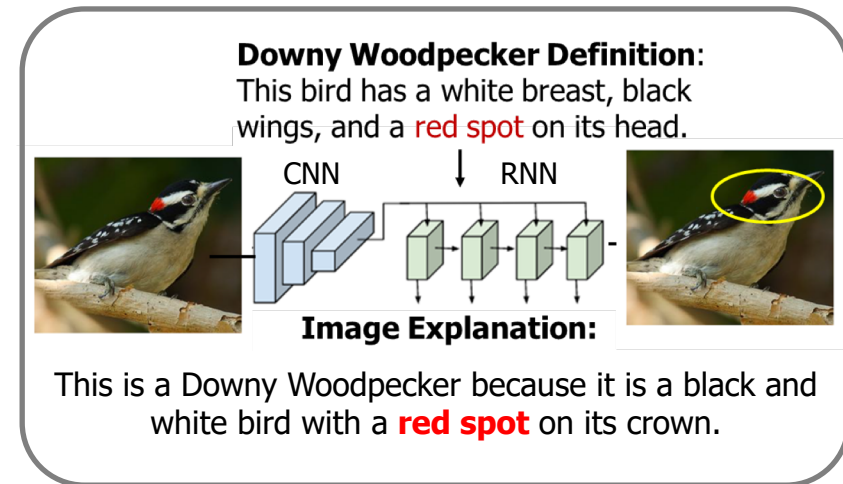| CP | Performer | Explainable Model | Explanation Interface |
|---|---|---|---|
| Both | UC Berkeley | Deep Learning | Reflexive and Rational |
| | Charles River | Causal Modeling | Narrative Generation |
| | UCLA | Pattern Theory+ | 3-level Explanation |
| Autonomy | Oregon State | Adaptive Programs | Acceptance Testing |
| | PARC | Cognitive Modeling | Interactive Training |
| | CMU | Explainable RL (XRL) | XRL Interaction |
| Analytics | SRI International | Deep Learning | Show and Tell Explanation |
| | Raytheon BBN | Deep Learning | Argumentation and Pedagogy |
| | UT Dallas | Probabilistic Logic | Decision Diagrams |
| | Texas A&M | Mimic Learning | Interactive Visualization |
| | Rutgers | Model Induction | Bayesian Teaching |

# Approaches to Deep Explanation

## Attention Mechanisms



Input Image   Salience Map

## Modular Networks



**Neural module networks**
[Andreas et al.CVPR16,EMNLP16] [Hu et al. CVPR17]

Q: Can you park here?

NO   Prediction

Neural module network

Attention visualization

Decision path

## Feature Identification



Generator

Target network

Tagger → Interpretation

## Learn to Explain



**Downy Woodpecker Definition**: This bird has a white breast, black wings, and a red spot on its head.

CNN    RNN

**Image Explanation:**

This is a Downy Woodpecker because it is a black and white bird with a **red spot** on its crown.

# Deeply Explainable Artificial Intelligence

## UC Berkeley, Boston U., U. Amsterdam, Kitware

| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|

### Deep Learning
- Post-hoc explanations by training additional DL models
- Explicit introspective explanations (Neural Module Networks)
- Reinforcement Learning
  - Informative rollouts
  - Explicit modular agent

### Reflexive and Rational
- Reflexive explanations (arise from the model)
- Rational explanations (come from reasoning about user's beliefs)
- Evaluation criteria
  - Human interpretability
  - Predictive behavior
  - Appropriate trust

### Autonomy
- Vehicle control (BDD-X, CARLA)
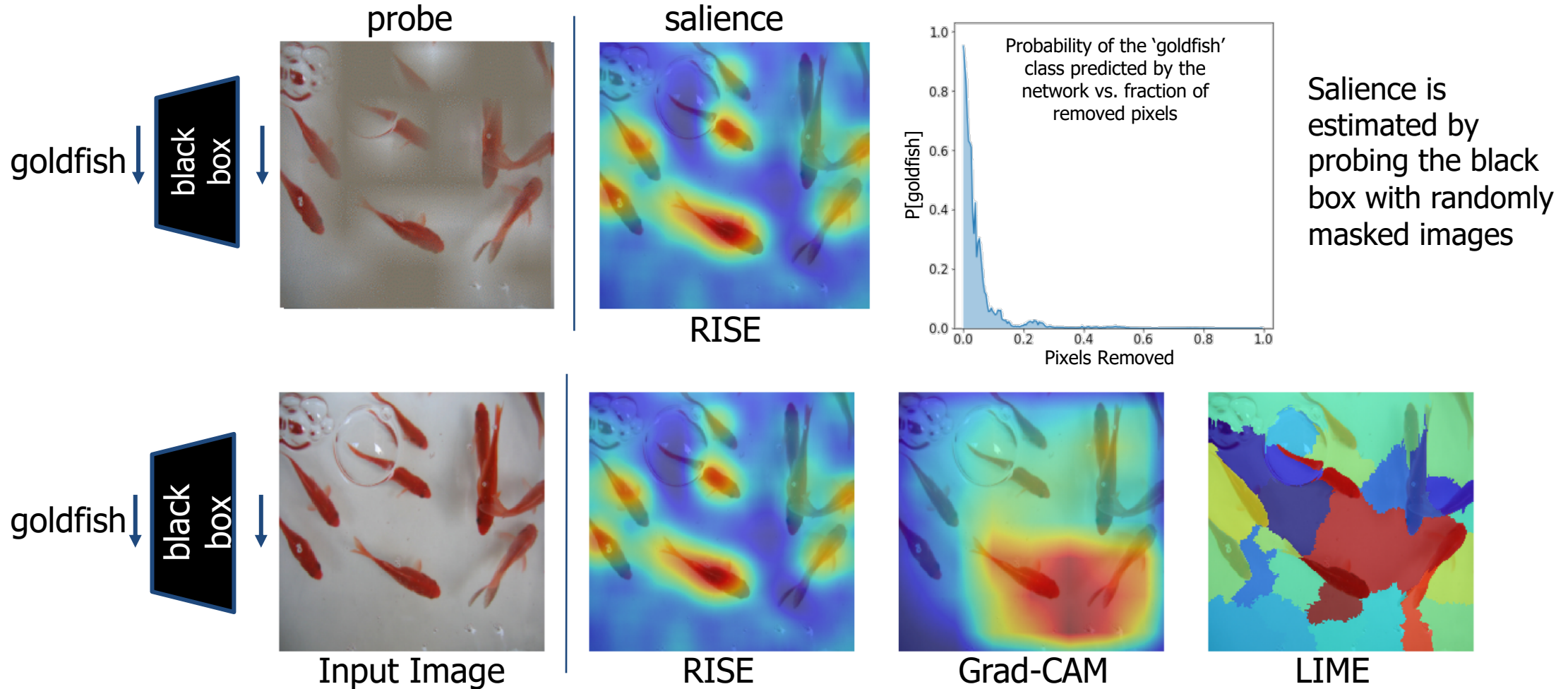- Strategy games (StarCraft II)

### Data Analytics
- Visual QA and filtering tasks (VQA-X, ACT-X, xView, DiDeMo, etc.)

---

- **PI**: Trevor Darrell (UC Berkeley)

---

- Pieter Abbeel (UC Berkeley)
- Tom Griffiths (UC Berkeley)
- Kate Saenko (Boston U.)
- Zeynep Akata (U. Amsterdam)
- Dan Klein (UC Berkeley)
- John Canny (UC Berkeley)
- Anca Dragan (UC Berkeley)
- Anthony Hoogs (Kitware)

# High Fidelity Visual Salience Model

## UC Berkeley, Boston U., U. Amsterdam, Kitware



probe

salience

goldfish

black box

RISE

Probability of the 'goldfish' class predicted by the network vs. fraction of removed pixels

Pixels Removed

Salience is estimated by probing the black box with randomly masked images

goldfish

black box

Input Image

RISE

Grad-CAM

LIME

Petsiuk, Das and Saenko. RISE: Randomized Input Sampling for Explanation of Black-box Models, 2018

Given the multi-modal explanation generated by the model, do you think the system will answer correctly?

Question: *Does this elephant have tusks?*

*"because there are no bones sticking out from its mouth"*



**Yes** **No**

<u>Incorrect!</u> The system answered *"no"* when the ground-truth answer is *"yes"*

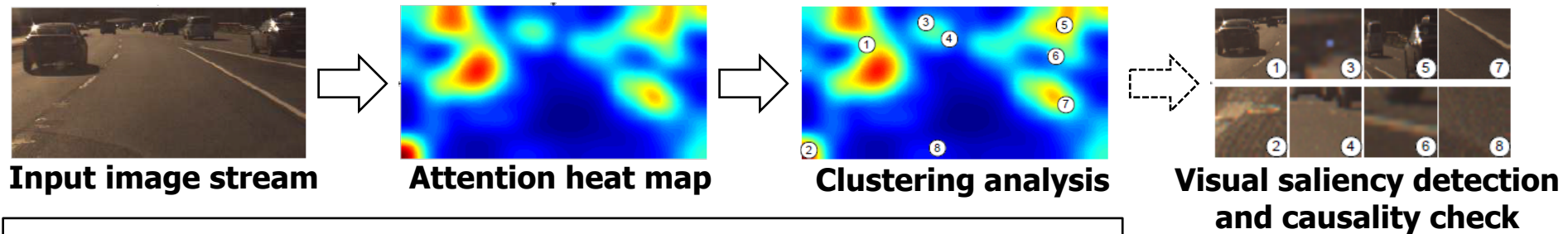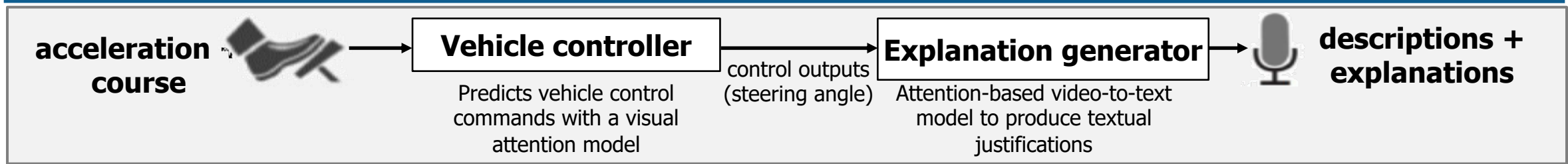Question: *Is this a professional sporting event?*

*"because the players are wearing official jerseys"*



**Yes** **No**

<u>Correct!</u> The system answered *"yes"* when the ground-truth answer is *"yes"*

| Explanation Effectiveness | Attention for Explanation Used? | Accuracy of Users Judgement |
|---|---|---|
| **Without explanation (existing SOTA)** | No | 57.5% |
| **UCB Model on descriptions** | Yes | 66.5% |
| **UCB Model without attention** | No | 61.5% |
| **UCB Model** | Yes | **70.0%** |

Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. Multimodal Explanations: Justifying Decisions and Pointing to the Evidence, 2018

# Causal Grounded Driving

acceleration course → **Vehicle controller** → control outputs (steering angle) → **Explanation generator** → descriptions + explanations

**Vehicle controller**
Predicts vehicle control commands with a visual attention model

**Explanation generator**
Attention-based video-to-text model to produce textual justifications



**Input image stream**

**Attention heat map**

**Clustering analysis**

**Visual saliency detection and causality check**

*Kim and Canny, in ICCV, 2017*
*Kim, Rohrbach, Darrell, Canny, and Akata, in NIPS Interpretable ML Symposium, 2017*

# CAMEL: Causal Models to Explain Learning

## Charles River Analytics (CRA), U. Mass, Brown

| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Causal Modeling**<br><br>• Experiment with the learned model (as a grey box) to learn an explainable, causal, probabilistic programming model | **Narrative Generation**<br><br>• Interactive visualization based on the generation of temporal, spatial narratives from the causal, probabilistic models | **Autonomy**<br>• Atari<br>• Starcraft<br><br>**Data Analytics**<br>• Pedestrian Detection (INRIA)<br>• Activity Recognition (ActivityNet) |

• **PI**: James Tittle (CRA)

• Jeff Druce (CRA)
• Avi Pfeffer (CRA)
• David Jensen (U. Mass)
• Michael Littman (Brown U.)

• James Niehaus (CRA)
• Emilie Roth (Roth Cognitive Engineering)
• Joe Gorman(CRA)
• James Tittle (CRA)

## Charles River Analytics, U. Mass, Brown

Generate causal explanations of ML operation and present them to the user as intuitive narratives in an interactive, easy-to-use interface grounded in cognitive engineering theories

**UCLA, Oregon State, Michigan State**

## Explainable Model

### Pattern Theory+

Interpretable representations
- STC-AOG: spatial, temporal, and causal models
- STC-PG: scene and event interpretations in analytics
- STC-PG+: task plans in autonomy

Theory of mind representations
- User's beliefs
- User's mental model of agent

## Explanation Interface

### 3-Level Explanation

- Concept compositions
- Causal and counterfactual reasoning
- Utility explanations

Explanation representations:
- X-AOG: explanation model
- X-PG: explanatory parse graph as dialogue
- X-Utility: priority and loss for explanations

## Challenge Problem

### Autonomy

- Robot executing daily tasks in physics-realistic VR platform
- Autonomous vehicle driving (GTA5 game engine)

### Data Analytics

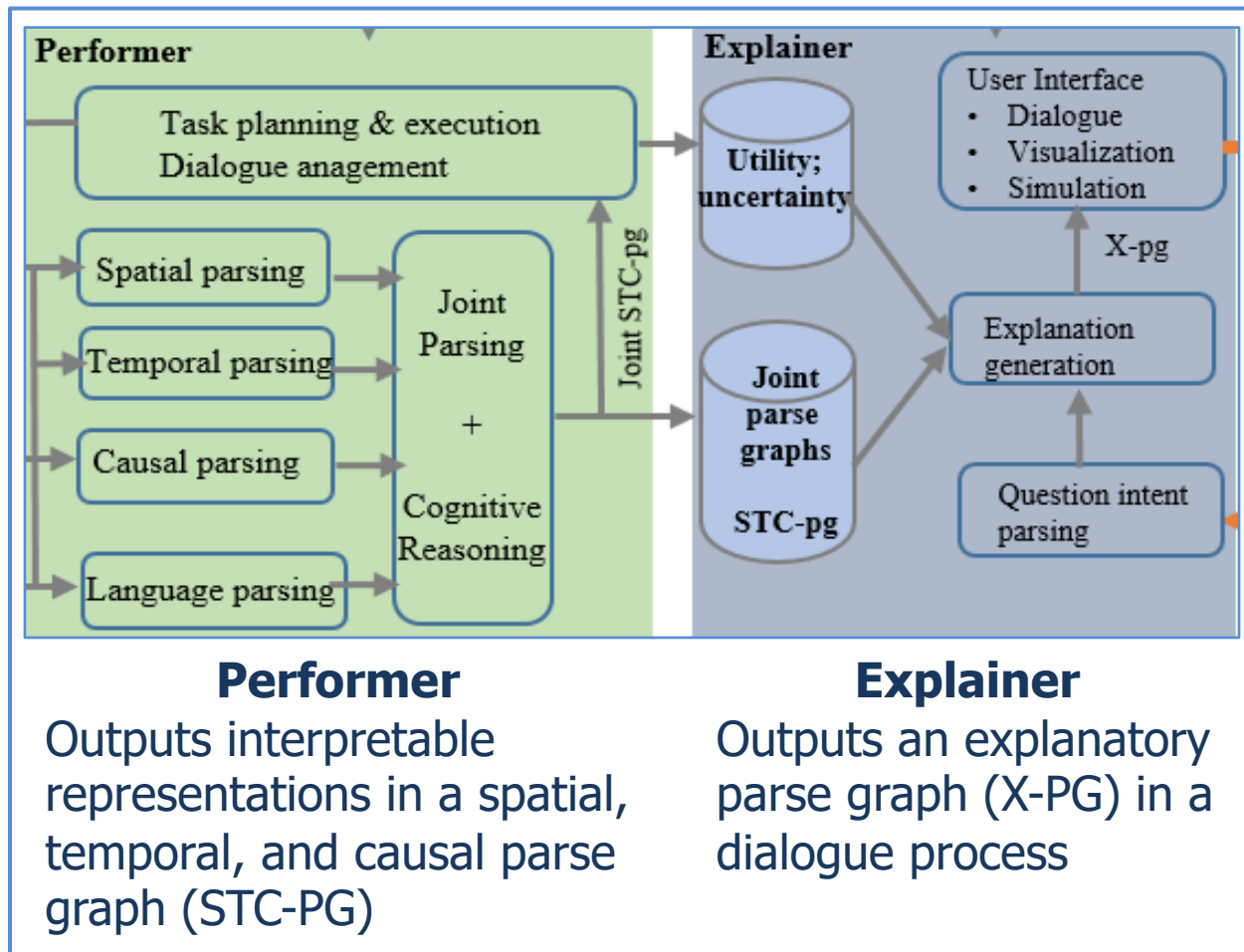- Network of video cameras for scene understanding and event analysis

- **PI**: Song-Chun Zhu (UCLA)

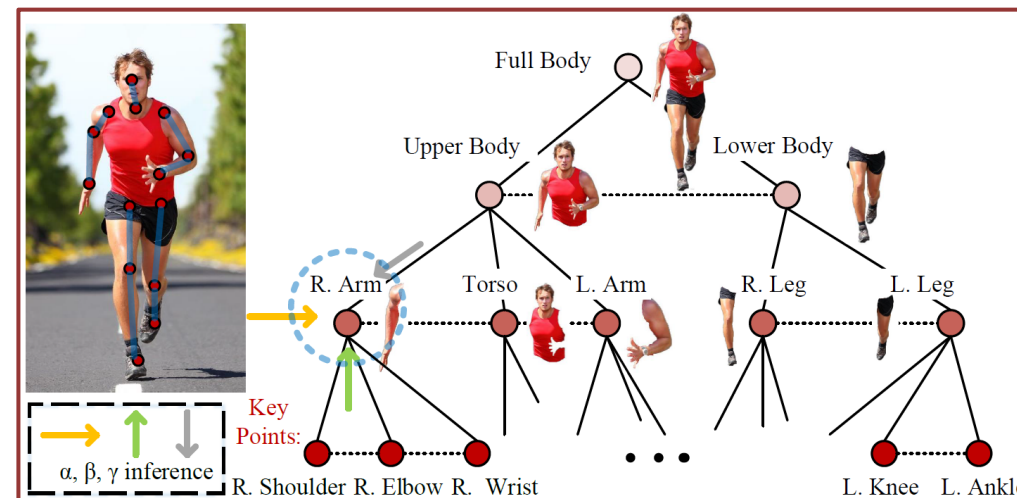- Ying Nian Wu (UCLA)
- Sinisa Todorovic (OSU)

- Joyce Chai (Michigan State)

## UCLA, Oregon State, Michigan State

### System Architecture



**Performer**
Outputs interpretable representations in a spatial, temporal, and causal parse graph (STC-PG)

**Explainer**
Outputs an explanatory parse graph (X-PG) in a dialogue process

### STC Parse Graph



An attributed parse graph for a running person. Each node has 3 computing channels:
- $\alpha$: grounding the node on DNN features;
- $\beta$: bottom-up;
- $\gamma$: top-down.

An explanation is represented as parse graph X-pg

# xACT: Explanation-Informed Acceptance Testing of Deep Adaptive Programs

## Oregon State University

### Explainable Model

**Adaptive Programs**

- Explainable Deep Adaptive Programs (xDAPs) – a new combination of Adaptive Programs, Deep Learning, and explainability

### Explanation Interface

**Acceptance Testing**

- Provides a visual and Natural Language explanation interface for acceptance testing by test pilots based on Information Foraging Theory
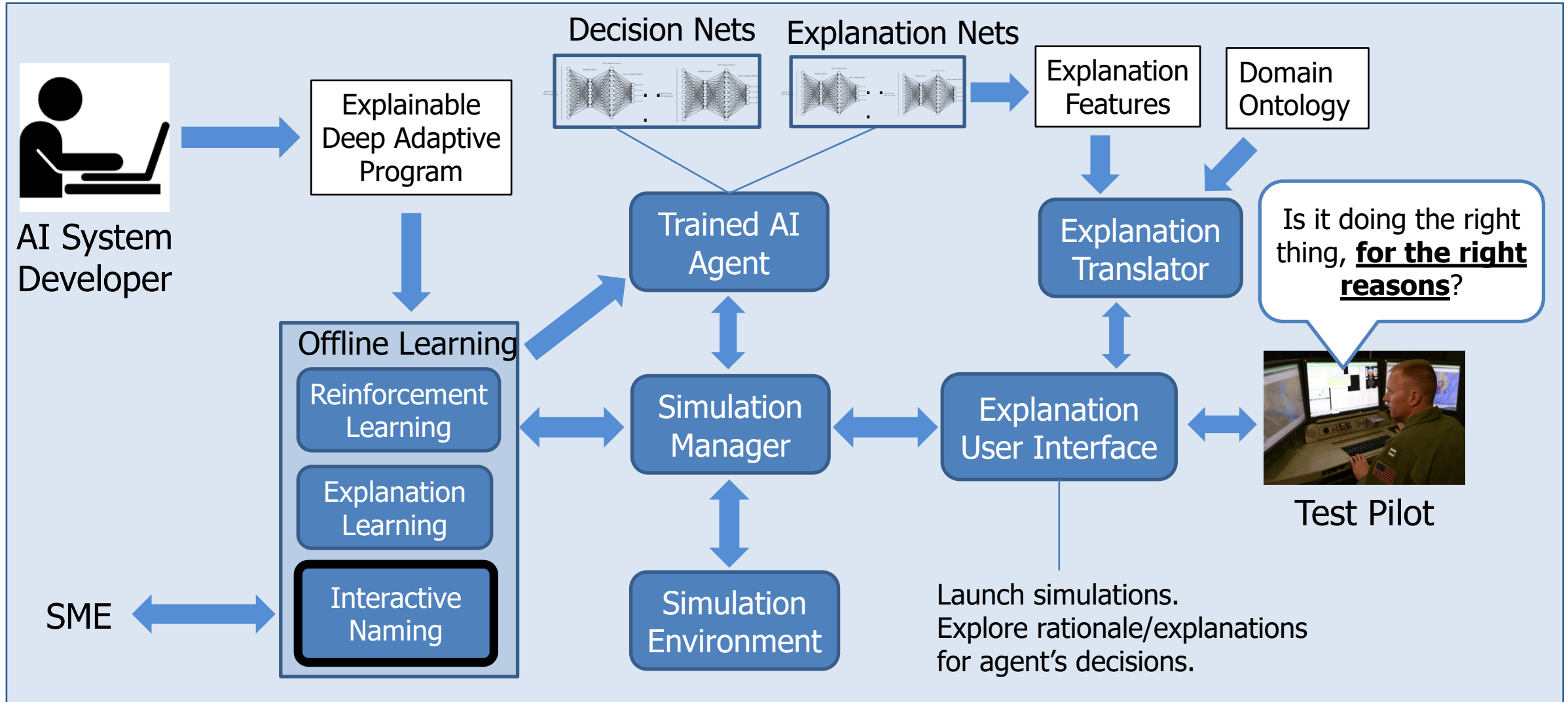
### Challenge Problem

**Autonomy**

- Real-time Strategy Games based on custom designed game engine designed to support explanation
- Starcraft II

- **PI**: Alan Fern (OSU)

- Tom Dietterich (OSU)
- Fuxin Li (OSU)
- Prasad Tadepalli (OSU)
- Weng-Keen Wong (OSU)

- Margaret Burnett (OSU)
- Martin Erwig (OSU)
- Liang Huang (OSU)

## Oregon State University

# COGLE: Common Ground Learning and Explanation

**PARC, CMU, U. Edinburgh, U. Michigan, USMA, IHMC**

## Explainable Model

### Cognitive Model

3-layer architecture
- Learning Layer (DNNs)
- Cognitive Layer (ACT-R Cognitive Model)
- Explanation Layer (HCI)

## Explanation Interface

### Interactive Training

- Interactive visualization of states, actions, policies, and values
- Module for test pilots to refine and train the system

## Challenge Problem

### Autonomy

- MAVSim wrapper over ArduPilot simulation environment
- Value of Explanation framework for measuring explanation effectiveness
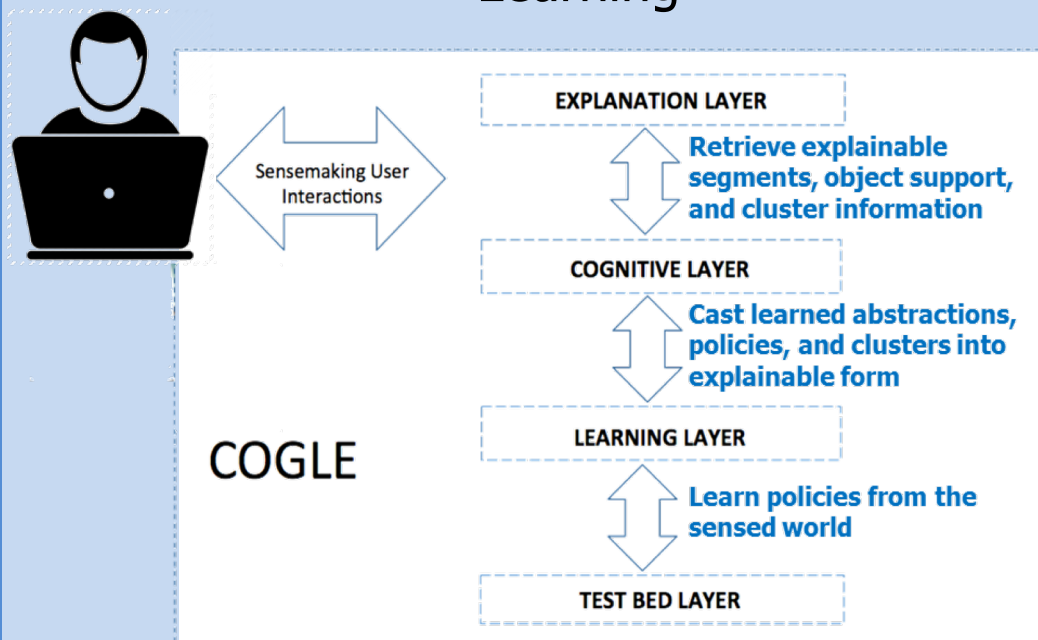
---

- **PI**: Mark Stefik (PARC)

---

- Honglak Lee (U. Michigan)
- Subramanian Ramamoorthy (U. Edinburgh)

- Christian Lebiere (CMU)
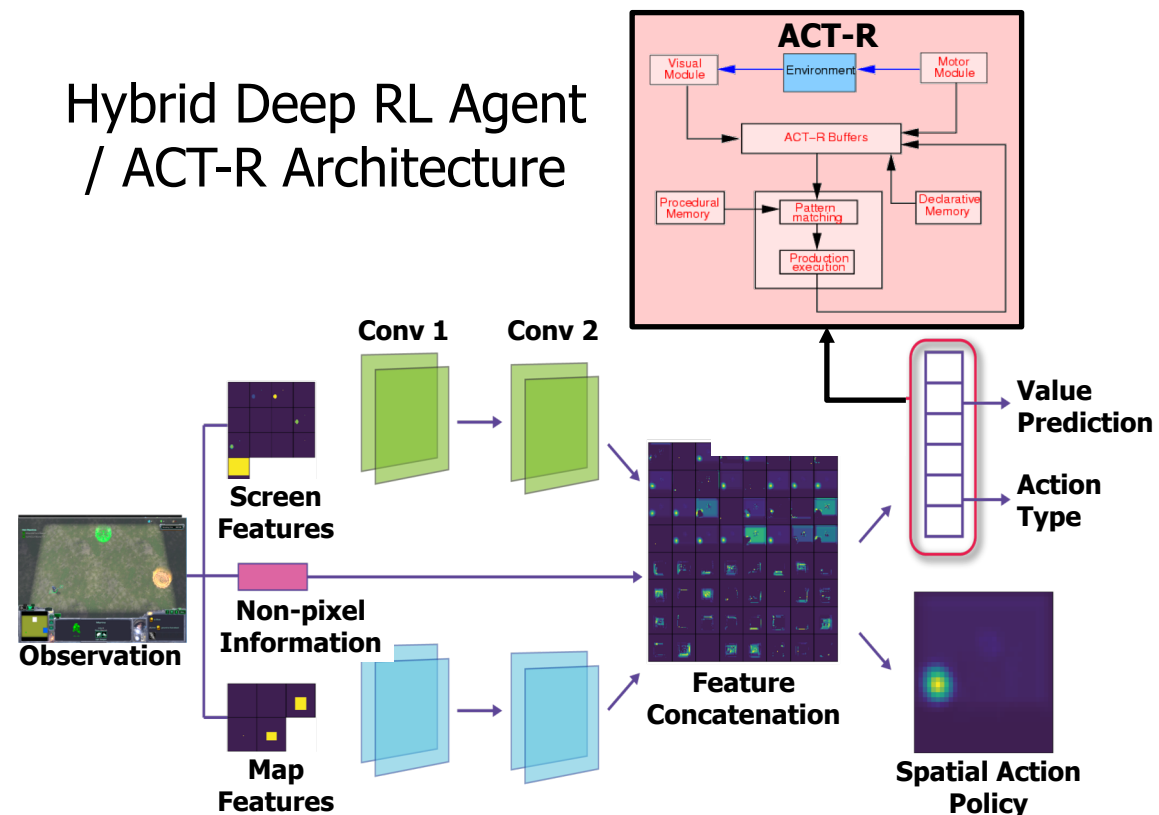- John Anderson (CMU)
- Robert Thomson (USMA)

- Michael Youngblood (PARC)

**PARC, CMU, U. Edinburgh, U. Michigan, USMA, IHMC**

## Layered Cognitive Architecture to Partition Explanation And Learning



## Hybrid Deep RL Agent / ACT-R Architecture

# XRL: Explainable Reinforcement Learning

## Carnegie Mellon University

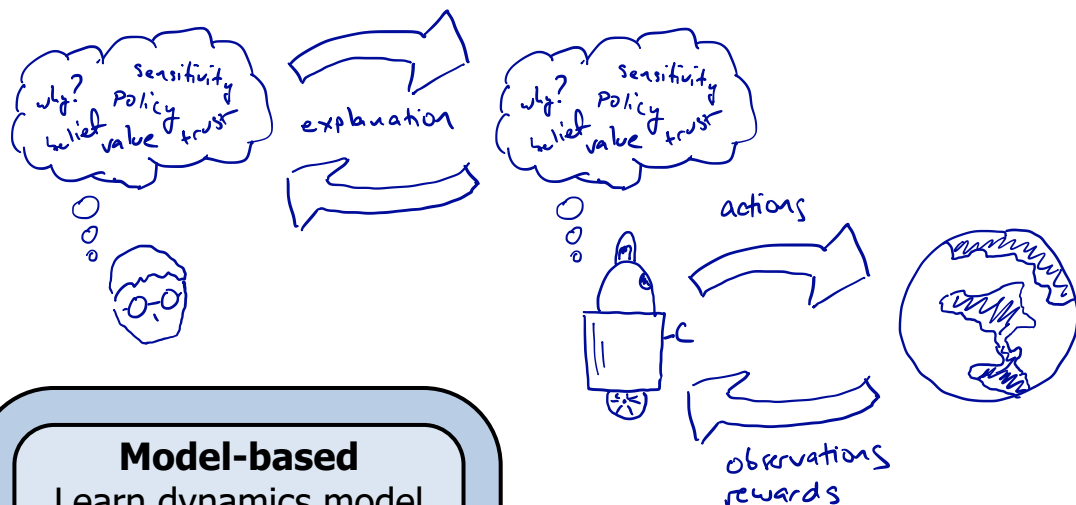| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Explainable RL (XRL)** | **XRL Interaction** | **Autonomy** |
| • Create a new scientific discipline for Explainable Reinforcement Learning with work on new algorithms and representations | • Interactive explanations of dynamic systems<br>• Human-machine interaction to improve performance | • Open AI Gym<br>• Autonomy in the electrical grid<br>• Mobile service robots<br>• Self-improving educational software |

• **PI**: Zico Kolter (CMU)

• Geoff Gordon (CMU)
• Pradeep Ravikumar (CMU)

## Carnegie Mellon University

Create a new discipline of explainable RL to enable dynamic human-machine interaction and adaptation for maximum team performance
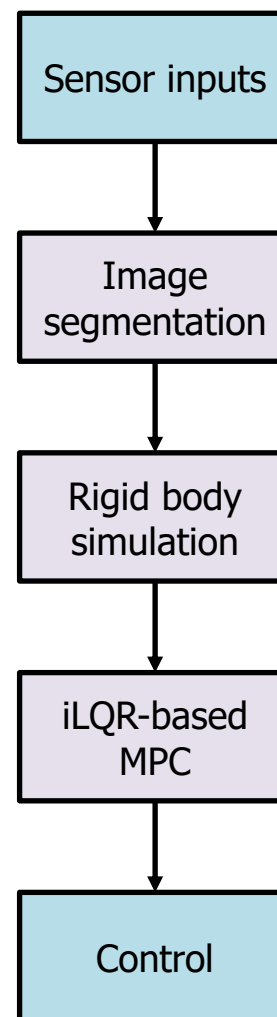


**Model-based**
Learn dynamics model of environment, plan actions in model, and execute in real system

Improve model learning/ creation for RL agents to capture benefits of model-based approach

**Model-free**
Directly learn value and/or policy for the environment

For any type of RL approach, provide an explanation of why an agent acted in a certain way

Sensor inputs

Image segmentation

Rigid body simulation

iLQR-based MPC

Control



**Differentiable Physics -** Applies implicit differentiation to solutions of LCP to analytically derive a backpropagation update of next state with respect to previous state, control, and model parameters

## SRI International, U. Toronto, UCSD, U. Guelph

| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Deep Learning** | **Show & Tell Explanation** | **Data Analytics** |
| Multiple deep learning techniques<br>• Attention-based mechanisms<br>• Compositional NMNs<br>• GANs | • DNN visualization<br>• Query evidence that explains DNN decisions<br>• Generate natural language justifications | • VQA<br>  • Visual Gnome<br>  • Flickr30<br>• MovieQA |

• **PIs**: Giedrius Burachas (SRI), Mohamed Amer (SRI)

| | | |
|---|---|---|
| • Xiao Lin (SRI) | Richard R. Zemel (U. Toronto) | • Jürgen Schulze (UCSD) |
| • Ryan Villamil (SRI) | Sanja Fidler (U. Toronto) | |
| • Dejan Jovanovic (SRI) | David Duvenaud (U. Toronto) | |
| • Avi Ziskind (SRI) | Graham Taylor (U. Guelph) | |
| • Michael Wessel (SRI) | | |

## SRI International, U. Toronto, UCSD, U. Guelph

Interpretable, Scene Graph-based VQA System with Active Attention



- Generate "show-and-tell" explanations with justifications of decisions accompanied by visualizations of input data used to generate inferences

- Scene and Situation Graphs, inferred from images and videos, support rich multimodal data analytics and explanations

- Scene Graphs guide attentional scanning for interpretable analytics

# EQUAS: Explainable QUestion Answering System

## Raytheon BBN, Georgia Tech, UT Austin, MIT

### Explainable Model

**Deep Learning**

- Semantic labelling of DNN neurons
- DNN audit trail construction
- Gradient-weighted Class Activation Mapping

### Explanation Interface

**Argumentation Theory**

- Comprehensive strategy based on argumentation theory
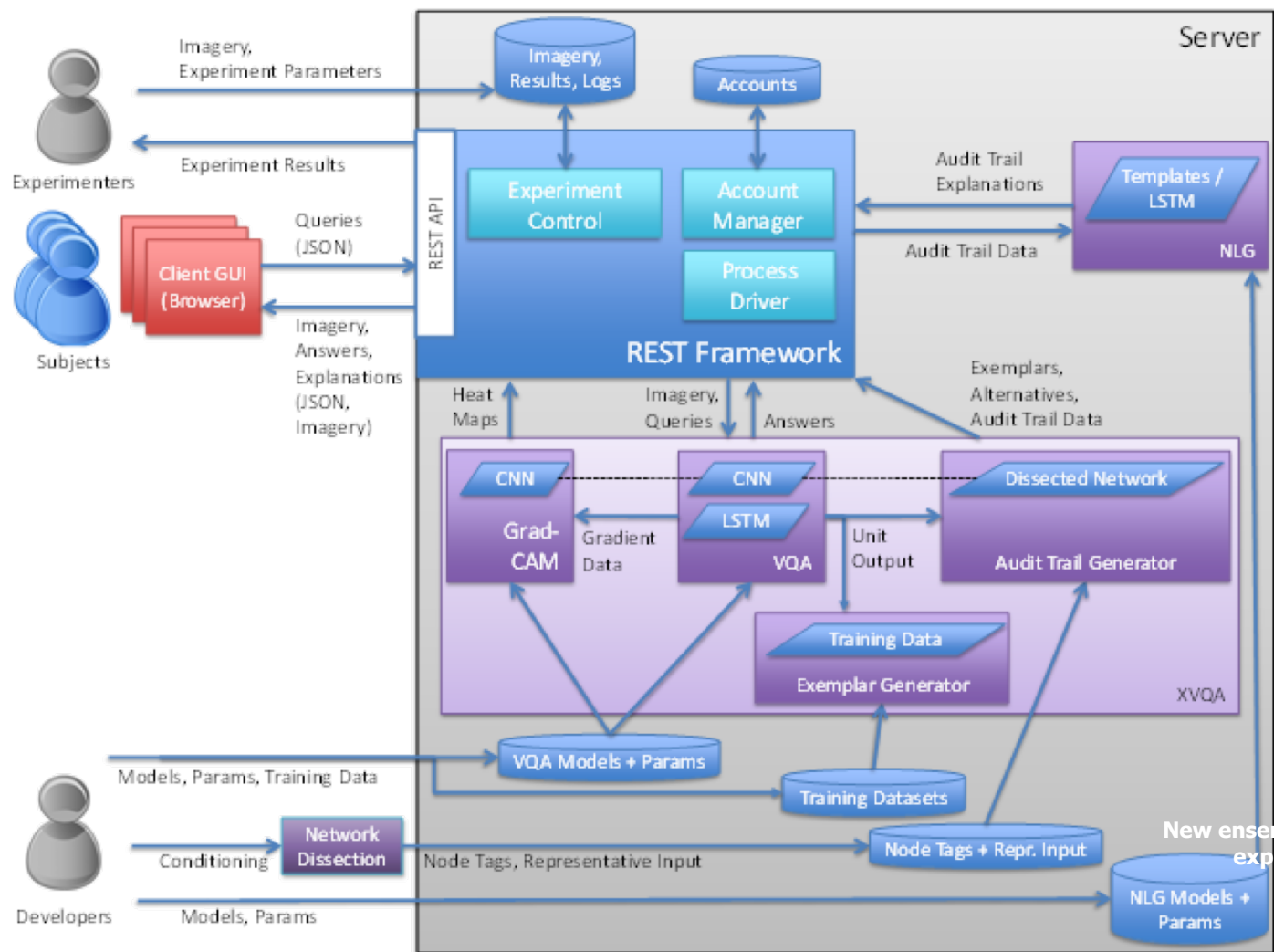- NL generation
- DNN visualization

### Challenge Problem

**Data Analytics**

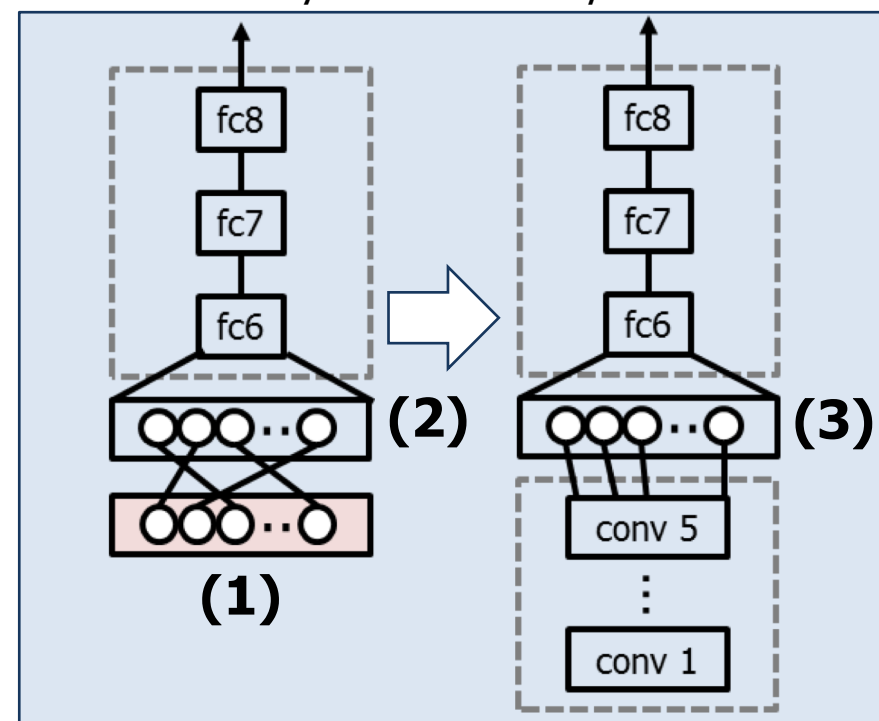- VQA for images and video

---

- **PI**: William Ferguson (Raytheon BBN)

---

- Antonio Torralba (MIT)
- Ray Mooney (UT Austin)
- Devi Parikh (Georgia Tech)
- Dhruv Batra (Georgia Tech)

## Raytheon BBN, Georgia Tech, UT Austin, MIT



Improve the interpretability of units using a **new conditioning method** to retrain the network to intentionally include *concept detectors*

1) **Pick units from standard vocabulary**
2) **Train top part of net**
3) **Use top to train bottom**

## UT Dallas, UCLA, Texas A&M, Indian Institute of Technology

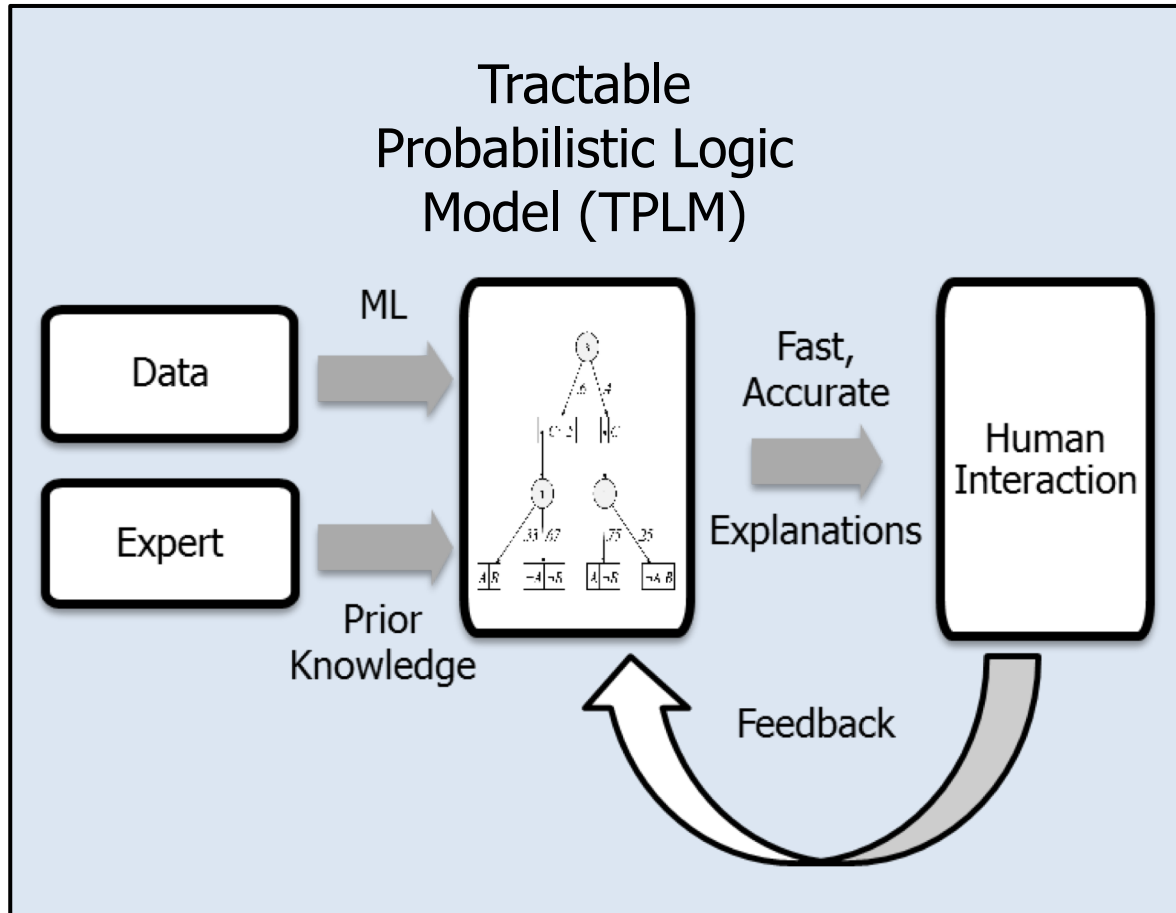| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Probabilistic Logic** | **Decision Diagrams** | **Data Analytics** |
| • Tractable Probabilistic Logic Models (TPLMs) – an important class of (non-deep learning) interpretable models | • Enables users to explore and correct the underlying model as well as add background knowledge | • Infer activities in multimodal data (video and text)<br>• Wetlab (biology) and TACoS (cooking) datasets |

- **PI**: Vibhav Gogate (UT Dallas)

- Adnan Darwiche (UCLA)
- Guy Van Den Broeck (UCLA)
- Nicholas Ruozzi (UT Dallas)
- Eric Ragan (Texas A&M)
- Parag Singla (IIT-Delhi)

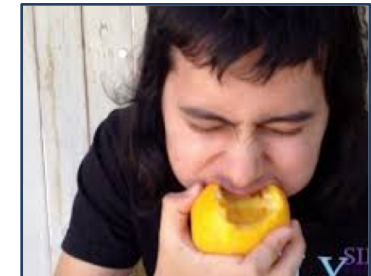## UT Dallas, UCLA, Texas A&M, Indian Institute of Technology

Use interpretable and tractable models based on well-founded principles from logic and probability theory



Tractable
Probabilistic Logic
Model (TPLM)

Data → ML → Fast, Accurate, Explanations → Human Interaction

Expert → Prior Knowledge

Feedback

**Find all videos in which a person peels oranges**
(explanations are captions generated by TPLMs)



Person using his hands to peel oranges. I can see the orange skin



Person using his hands. I can see the orange skin and skinless orange



Person using his hands to peel oranges. I can see the orange skin on the table and peeled oranges

# Transforming Deep Learning to Harness the Interpretability of Shallow Models

## Texas A&M, Washington State

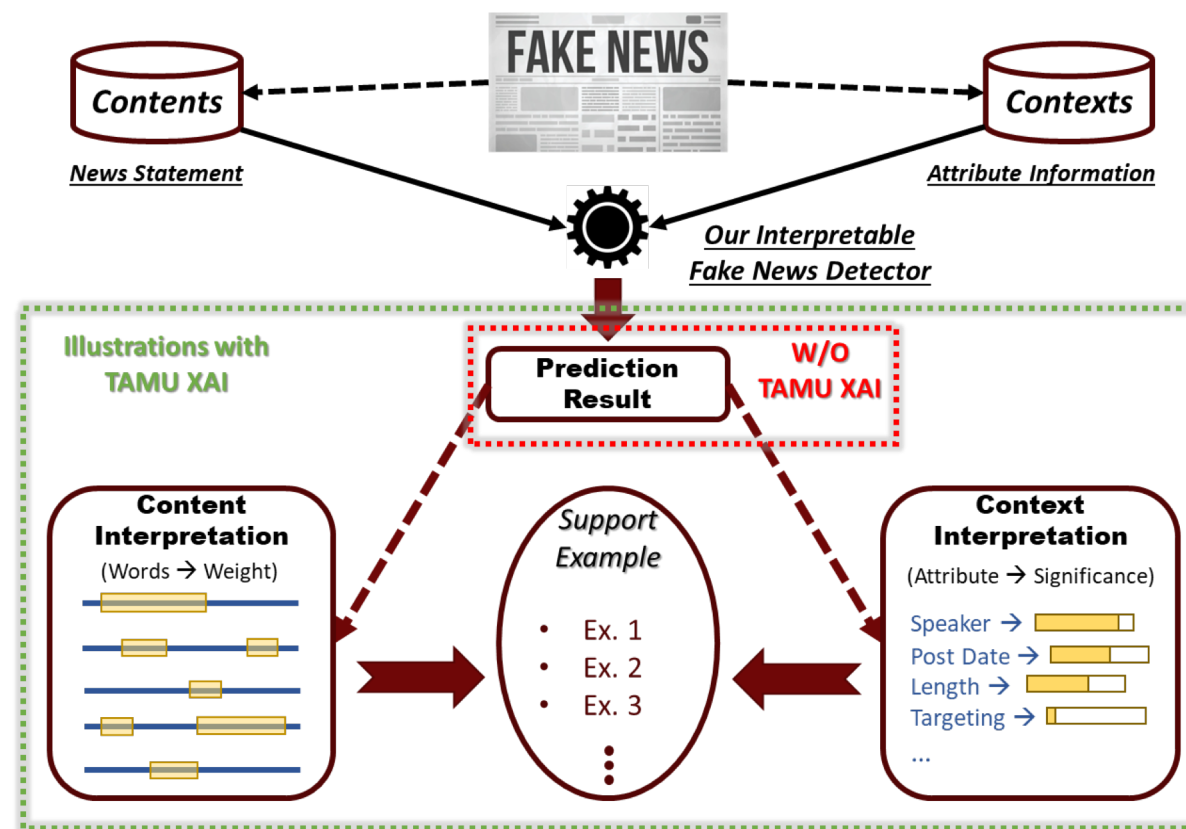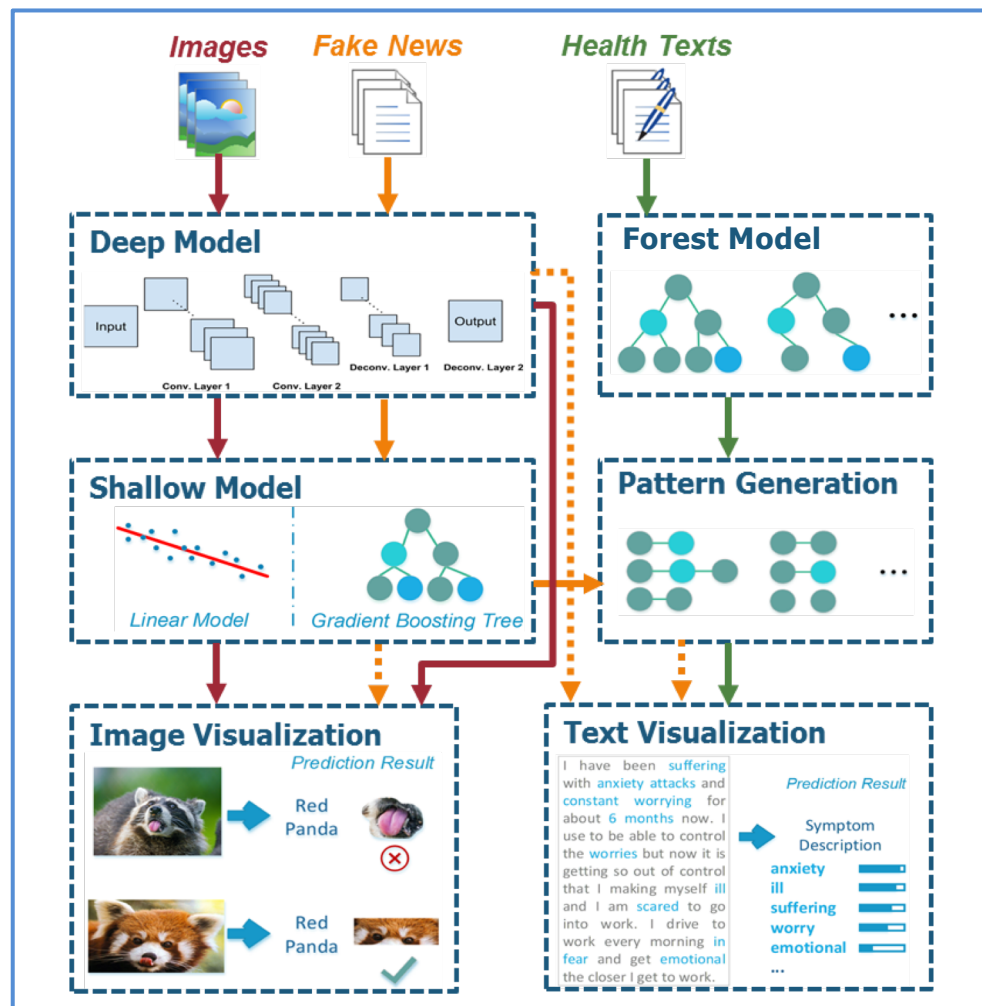| Explainable Model | Explanation Interface | Challenge Problem |
|---|---|---|
| **Mimic Learning** | **Interactive Visualization** | **Data Analytics** |
| • Mimic learning framework combines DL models for prediction and shallow models for explanations<br>• Interpretable learning algorithms extract knowledge from DNNs for relevant explanations | • Interactive visualization over multiple views, using heat maps and topic modeling clusters to show predictive features | • Multiple tasks using data from Twitter, Facebook, ImageNet, and news websites |

• **PI**: Xia Hu (Texas A&M)

• Shuiwang Ji (Wash. State)   • Eric Ragan (Texas A&M)

## Texas A&M, Washington State

Develop an end-to-end interpretable deep learning infrastructure with image and text datasets

# Model Explanation by Optimal Selection of Teaching Examples

## Rutgers University

### Explainable Model

**Model Induction**
- Select the optimal training examples to explain model decisions based on Bayesian Teaching

### Explanation Interface

**Bayesian Teaching**
- Example-based explanation of
  - Full model
  - User-selected sub-structure
  - User submitted examples

### Challenge Problem

**Data Analytics**
- Image processing
- Text corpora
- VQA
- Movie events

- **PI**: Patrick Shafto (Rutgers)

- Scott Cheng-Hsin Yang (Rutgers)

## Rutgers University

Extend Bayesian teaching to enable automatic explanation by selecting the subset of data that are most representative of the model's generative process



Good and bad examples for teaching a category (illustrates model strengths and weaknesses)

Good pairs of examples of the category 9



Bad pairs of examples of the category 9

# XAI Program Schedule



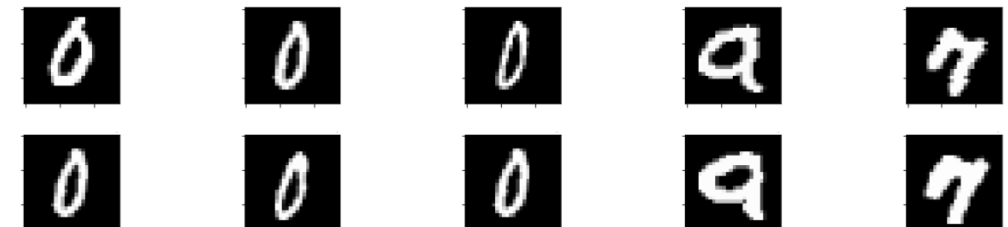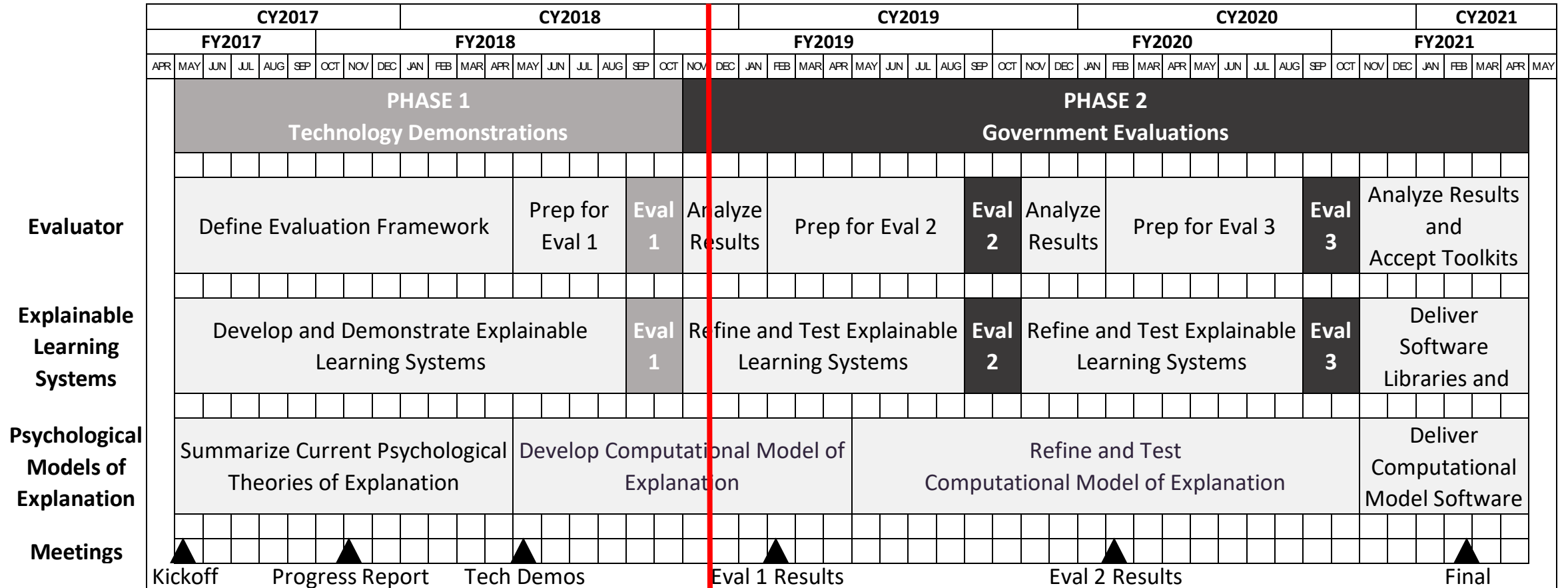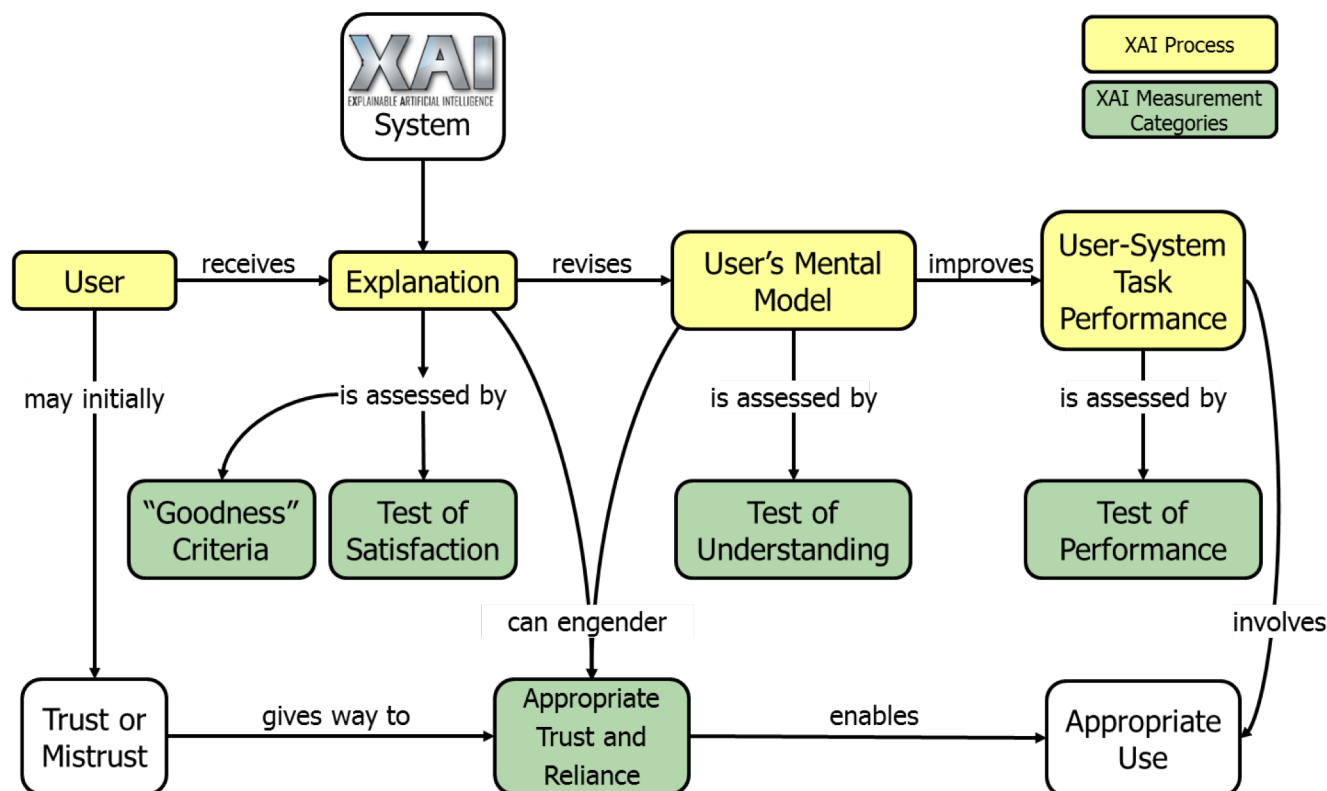| | CY2017 | | CY2018 | | | CY2019 | | CY2020 | | CY2021 |
|---|---|---|---|---|---|---|---|---|---|---|
| | FY2017 | FY2018 | | FY2019 | | FY2020 | | FY2021 | | |

Months row: APR MAY JUN JUL AUG SEP OCT NOV DEC JAN FEB MAR APR MAY JUN JUL AUG SEP OCT NOV DEC JAN FEB MAR APR MAY JUN JUL AUG SEP OCT NOV DEC JAN FEB MAR APR MAY JUN JUL AUG SEP OCT NOV DEC JAN FEB MAR APR MAY

**PHASE 1 — Technology Demonstrations**

**PHASE 2 — Government Evaluations**

## Evaluator
Define Evaluation Framework | Prep for Eval 1 | Eval 1 | Analyze Results | Prep for Eval 2 | Eval 2 | Analyze Results | Prep for Eval 3 | Eval 3 | Analyze Results and Accept Toolkits

## Explainable Learning Systems
Develop and Demonstrate Explainable Learning Systems | Eval 1 | Refine and Test Explainable Learning Systems | Eval 2 | Refine and Test Explainable Learning Systems | Eval 3 | Deliver Software Libraries and

## Psychological Models of Explanation
Summarize Current Psychological Theories of Explanation | Develop Computational Model of Explanation | Refine and Test Computational Model of Explanation | Deliver Computational Model Software

## Meetings
Kickoff | Progress Report | Tech Demos | Eval 1 Results | Eval 2 Results | Final

**Phase 1 Evaluations**

## Explanation Process & Measures



## Experimental Conditions

***Without Explanation*** - The explainable learning system is used to perform a task without providing an explanation to the user

***With Explanation*** - The explainable learning system is used to perform a task and generates explanations for every recommendation or decision it makes, and every action it takes

***Partial Explanation*** - The explainable learning system is used to perform a task and generates only partial or ablated explanations (to assess various explanation features)

***Control*** - A baseline state-of-the-art non-explainable system is used to perform a task

www.darpa.mil